



کاربرد فناوری ریز آرایه (Microarray) در تشخیص عوامل بیماری های عفونی و فرایند آنالیز داده های ریز آرایه

نظری خدیجه^{۱***}، دکتر کرمی علی^۲، دکتر مهدوی امیری نظام الدین^۳

(۱) کارشناس ارشد علوم کامپیوتر، دانشگاه صنعتی شریف (۲) دانشیار مرکز تحقیقات بیولوژی مولکولی، دانشگاه علوم پزشکی بقیه ا. (عج) (۳)
استاد دانشکده ریاضی، دانشگاه صنعتی شریف

چکیده

فناوری ریز آرایه یا میکروآری تلفیقی از توانمندی های علوم مختلف مانند بیولوژی مولکولی، میکروالکترونیک، میکروفلئیدیک و بیوانفورماتیک است. با استفاده از این فناوری می توان همان زمان هزاران هدف ژنتیکی و یا پروتئینی را در یک مجموعه ای کوچک بررسی کرد.

ریز آرایه های DNA عبارتند از مجموعه ای از نقاط میکروسکوپی DNA که به سطح جامدی مثل شیشه، پلاستیک یا تراشه هی سیلیکون متصل شده اند و یک آرایه را تشکیل می دهند. قطعات DNA ثابت شده به عنوان کاوشگر در نظر گرفته می شوند. در یک آزمایش می توان از هزاران کاوشگر استفاده کرد. بنابراین هر ریز آرایه شامل همین تعداد آزمون ژنتیکی است که برای همه ای آنها آزمایش به صورت موازی انجام می شود. با این توانایی ریز آرایه ها به بسیاری از بررسی های زیست شناختی سرعت بخشیده اند.

مهم ترین مرحله در این فناوری، تجزیه و تحلیل داده های انبوهی است که حاصل می شود و به ابزاری بیوانفورماتیکی که بتواند این داده ها را با درجه ای اطمینان بالا تحلیل کند، نیاز است.

بیماری های عفونی از آغاز زندگی بشری همواره همراه انسان بوده و باعث مشکلات اساسی شده اند. یکی از کاربردهای مهم فناوری ریز آرایه، امکان آزمودن وجود هزاران میکروارگانیسم در نمونه های محیطی یا بالینی تنها در یک آزمایش و در زمان کوتاهی است که منجر به تشخیص به موقع عامل عفونی و بیماری زا می شود و به این ترتیب گام مهمی در درمان بیماری برداشته ایم. همان طور که ذکر شد، مهم ترین مرحله در فناوری ریز آرایه تجزیه و تحلیل داده هاست. ما در اینجا الگوریتم E-Predict و بسته ای نرم افزاری DetectiV را که بر پایه ای تعیین گونه های پاتوژن در ریز آرایه است، تشریح می کنیم. در ادامه، کاربرد این دو روش در کشف عامل بیماری زا در مجموعه داده بزرگی که به صورت عمومی در دسترس است، بیان می کنیم و نشان می دهیم که DetectiV بهتر از E-Predict عمل می کند.

واژه های کلیدی : ۱- ریز آرایه ۲- داده های ریز آرایه ۳- تحلیل داده های ریز آرایه ۴- بیماری های عفونی

مقدمه

روش‌های قدیمی بررسی مطالعه و تحلیل بیان ژن از جمله RT-PCR، نورترن بلات و ساترن بلات و تعیین ردیف ژن از بهترین روش‌ها برای آنالیز تعداد محدودی ژن و نمونه در یک زمان است.^(۱) اما فناوری ریزآرایه قادر است همزمان هزاران ژن یا پروتئین را مورد بررسی قرار دهد. نوع ریزآرایه بستگی به موادی دارد که بر روی اسلايد قرار می‌گیرند: اگر DNA باشد به آن ریزآرایه‌ی RNA، اگر RNA باشد ریزآرایه‌ی DNA. اگر پروتئین باشد، ریزآرایه‌ی پروتئین و اگر ماده‌ی به دست آمده از بافت ویژه‌ای باشد، ریزآرایه‌ی بافتی نامیده می‌شود.

مراحل آزمایشگاهی فناوری ریزآرایه شامل ساخت ریزآرایه، استخراج و نشان‌دار کردن نمونه با استفاده از رنگ‌های فلورسانسی، هیبرید کردن نمونه با ریزآرایه و شستشوی ریزآرایه است. در نهایت با استفاده از اسکنر شدت نور پروب‌های موجود در ریزآرایه اندازه‌گیری و کمی می‌شود و نوبت به مرحله‌ی تحلیل داده‌ها می‌رسد.^(۲)

ریزآرایه‌های DNA از ابزارهای مفید و نوظهور در متاثرnomیک هستند. این فناوری در آشکارسازی عوامل عفونی مختلف شامل انواع باکتری‌ها، ویروس‌ها، انگل‌ها، قارچ‌ها و همچنین کشف پاتوژن‌های انسانی نوظهور، بسیار موفق عمل کرده است.^(۳-۷) شاید مهم ترین کاربرد فناوری ریزآرایه در بیماری‌های عفونی، توانایی تشخیص سریع و همزمان چند هدف و همچنین تفکیک ژنتوتایپ پاتوژن‌ها می‌باشد است. به دلیل آن که پاتوژن‌ها ترکیبات ژنتیکی مجزا دارند و ریزآرایه‌ها قادرند که توالی همه‌ی ژن‌ها را همزمان بررسی کنند، این فناوری یک ابزار ایده‌آل برای این هدف است.^(۸,۹)

بسیاری از ریزآرایه‌ها با هدف مشخص کردن طیفی از میکروارگانیسم‌های موجود در یک نمونه طراحی شده اند و عملاً کاربرد آن‌ها در آشکارسازی ویروس‌های بیماری‌زا تایید شده است.^(۱۰) شناسایی انواع عوامل تهدید بیوتوریستی^(۱۱) و بررسی کیفیت آب و مواد غذایی^(۱۲) جهت تشخیص انواع عوامل بیماری‌زا از دیگر کاربردهای ریزآرایه است.

به هر حال آنالیز داده‌های ریزآرایه در بیماری‌های عفونی از آن جهت که هر نمونه‌ی اسیدنوکلئیکی مورد آزمون معمولاً شامل مخلوطی از DNA و RNA از ارگانیسم‌های مختلف، میزان و آلودگی‌های مختلف است و همه‌ی این‌ها ممکن است در نتیجه‌ی آزمایش اثرگذار باشند، کار آسانی نیست. علاوه بر این امکان دارد که با وجود چندین گونه‌ی پاتوژن و حتی به هم وابسته، عمل هیبریداسیون پیچیده‌تر شود.

مواد و روش‌ها

در این بررسی دو روش DetectiV و E-predict جهت تحلیل داده‌های ریزآرایه‌ی هدف، استفاده شده است.

E-predict

یکی از روش‌های تحلیل داده‌های ریزآرایه استفاده از الگوریتم E-Predict^(۱۳) است. اساس این روش به منظور شناسایی گونه‌ها، الگوهای هیبریداسیون مشاهده شده‌ی ریزآرایه است. در این الگوریتم، ابتدا پروفایل‌های انژی هیبریداسیون برای هر ژنوم ویروس مرجع به صورت کامل توالی یابی شده موجود در GenBank (شامل ۱۲۲۹ ویروس مجزا) محاسبه می‌شود.

انتظار می‌رود که همه‌ی الیگونوکلئوتیدهای ریزآرایه با ژنوم ویروس‌های معین که با استفاده از^۱BLAST^(۱۴) مشخص شده‌اند، هیبرید شوند. سپس انژی آزاد هیبریداسیون (ΔG) برای هر همترازی با استفاده از روش نزدیکترین همسایه محاسبه می‌شود. برای الیگوهایی که در تولید خروجی BLAST ناموفق هستند، انژی هیبریداسیون صفر در نظر گرفته می‌شود. بنابراین یک پروفایل انژی نظری مشخص شامل انژی‌های هیبریداسیون غیر صفری است که برای زیرمجموعه‌ی ای از الیگوها که یک همترازی BLAST مطابق با ژنوم تولید می‌کنند، محاسبه می‌شود. جمعاً، پروفایل‌های انژی از همه‌ی ویروس‌ها تشکیل

¹NCBI BLAST [http://www.ncbi.nlm.nih.gov/BLAST]

یک ماتریس انرژی می‌دهند. در این ماتریس، هرسطر مطابق با یک گونه‌ی ویروس و هر ستون مطابق با یک الیگو از یک زیرآرایه است.

در مرحله‌ی بعد بردار شدت الیگوها نرمال می‌شود و با هر پروفایل نرمال شده در ماتریس انرژی با استفاده از یک معیار مشابهت مقایسه می‌شود. نتیجه‌ی این عملکرد یک بردار از امتیازات مشابهت و به صورت خام است. هر عنصر در این بردار، مشابهت بین الگوهای مشاهده شده و یکی از پروفایل‌های پیش‌بینی شده برای یک گونه‌ی موجود در ماتریس انرژی را بیان می‌کند.

در روش E-Predict مقدار p همراه با امتیاز مشابهت جواب را مشخص می‌کند. با استفاده از E-Predict در تعدادی از وضعیت‌ها نتایج مفیدی حاصل شده است. در هر صورت، در حال حاضر E-Predict شامل ابزاری برای تصویرسازی نیست و احتیاج به بهنگام سازی و محاسبات پرهزینه قبل از استفاده‌ی آن برای آرایه‌های جدید دارد. همچنین E-Predict تنها در سیستم عامل Unix و Linux قابل استفاده است.

DetectiV

DetectiV به عنوان بسته‌ای برای نرم افزار آماری R شامل توابعی برای تصویرسازی، نرمال سازی و آزمون معنی‌داری از داده‌های ریزآرایه‌ی تشخیص پاتوژن است. R یک نرم افزار آماری رایگان است و در سیستم عامل‌های Windows, Unix, Mac قابل استفاده است.^(۱۵) از آن جا که DetectiV در نرم افزار R ایجاد می‌شود، به آسانی با بسیاری از بسته‌هایی که برای تحلیل ریزآرایه در دسترس هستند، ادغام می‌شود.^(۱۶)

مجموعه داده‌ی اصلی یک ماتریس از داده‌های است که سطرهای آن نمایش دهنده‌ی پروب‌ها و ستون‌ها نمایش دهنده‌ی اندازه‌گیری‌ها از ریزآرایه هاست. این داده‌ها به آسانی از ساختار داده‌هایی که به وسیله‌ی limma به دست آمده‌اند، ایجاد می‌شوند.^(۱۷) DetectiV شامل توابعی برای خواندن فرمتهای خروجی اکثر اسکنرهای affy داری توابعی برای خواندن داده‌های افی متريکس است.^(۱۸)

روی پروب‌های تکراری ميانگين می‌گيريم و در مرحله‌ی بعدی عمل نرمال‌سازی را با استفاده از تقسيم مقادير متناظر با هر پروب بر مقدار ميانه‌ی آرایه‌ی متناظر و سپس لگاريتم گيری در مبنای ۲ انجام می‌دهيم. به اين ترتيب داده‌ها نرمال و ميانگين آن‌ها صفر است. (البته فرض بر اين است که بيشتر پروب‌ها با چيزی هيبريد نمي‌شوند).

بعد از آن با استفاده از آزمون t داده‌ها براساس گونه‌ی ویروس گروه‌بندی می‌شوند. داده‌هایی در مراحل بعدی استفاده می‌شوند که لگاريتم ميانگين در مبنای ۲ آن‌ها، بزرگتر یا مساوی ۱ باشد. اين اعداد بر اساس مقدار p مرتب می‌شوند. در اين مبحث داده‌های استفاده شده از GEO با عدد دست‌يابي GSE2228 پياده شده‌است.^(۱۹) پلت فرم آرایه برای اين داده‌ها، GEO^۱ با عدد دست‌يابي GPL1834 است و شامل بيش از ۱۱۰۰۰ الیگوست که نشان دهنده‌ی ۱۰۰۰ گونه‌ی ویروس و باكتري است.

نتایج و بحث

نتایج مورد انتظار و صحیح از قبل مشخص شده‌اند در صورتی نتایج به دست آمده از تحلیل داده‌ها صحیح فرض می‌شود که یکی از دو مورد زیر اتفاق بیافتد:

- (۱) جواب حاصل از تحلیل با عامل بیماریزای از قبل مشخص مشابه باشد.
- (۲) اگر در آرایه عامل بیماریزای موردنظر وجود نداشته باشد، آنگاه نتیجه‌ی حاصل بسیار مشابه و وابسته به ویروس مورد نظر باشد.

^۱ Gene Expression Omnibus[<http://ncbi.nlm.nih.gov/GEO>



در ۵۵ آرایه از ۵۶ آرایه موجود منجر به نتیجه‌ی صحیح می‌شود. در مقایسه، روش DetectiV در ۵۳ آرایه از ۵۶ آرایه نتیجه‌ی صحیح می‌دهد.

DetectiV از تکنیک‌های خوب و ساده‌ی تصویرسازی و روش‌های آماری مناسب برای تحلیل داده‌های ریزآرایه‌های آشکارسازی پاتوژن استفاده می‌کند.

ما در اینجا بین مفاهیم آماری و زیستی تفاوت قابل می‌شویم. در روش DetectiV کمترین مقدار p در ترکیب با بزرگترین مقدار میانگین، و در روش E-Predict کمترین مقدار p همراه با بزرگترین مقدار مشابه‌ی جواب را مشخص می‌کند. در بسیاری از حالات با استفاده از قوانین خودکار می‌توانیم جواب درست را نتیجه بگیریم. به هر حال، به ناچار بعضی حالات وجود دارند که نیازمند دخالت انسان است.

نتایج کاربرد DetectiV روی مجموعه داده‌ی GSE8746 با جواب درست در هر ۱۲ آرایه مشوق بسیار خوبی برای ماست. بهویژه، توانایی آرایه و DetectiV در تمیز دادن بین نه تنها گونه‌های ویروس بلکه زیرنوع‌های FMDV بسیار جذاب و درخور توجه است و نمایانگر قدرت آن است.

نتایج کاربرد DetectiV روی مجموعه داده‌ی سارس نیز مشوق و امیدوار کننده است. در این‌جا، الیگوهایی که مخصوص ویروس سارس طراحی شده باشند، روی آرایه وجود ندارند.

این نکته شایان ذکر است که برای کاربرد DetectiV در مجموعه داده‌ی دیگری که از آرایه‌های کاملاً متفاوت با مجموعه داده‌ی اول استفاده می‌کند، تنها نیاز است که کاربر عدد دست‌یابی GEO و تعداد آرایه‌های داخل مجموعه داده را تغییر دهد. این را می‌توان با E-Predict مقایسه کرد که نیازمند محاسبه‌ی ماتریس مشابه‌ی پیچیده و بزرگ و همچنین بهینه‌سازی چندین پارامتر است.

References

1. Brown T.A.N: Gene cloning and DNA analysis an introduction. *Blackwell Science 2006*.
2. Wentian L, Yaning Y : Introduction to microarray analysis. *BMC proceedings 2007*
3. Riesenfeld CS, Schloss PD, Handelsman J: Metagenomics: genomic analysis of microbial communities. *Annu Rev Genet 2004*.
4. Lapa S, Mikheev M, Shchelkunov S, Mikhailovich V, Sobolev A, Blinov V, Babkin I, Guskov A, Sokunova E, Zasedatelev A, et al.:Species level identification of orthopoxviruses with an oligonucleotide microchip. *J clin Microbiol 2002*.
5. Boonham N, Walsh K, Smith P, Madagan K, Graham I, Barker I: Detection of potato viruses using microarray technology:towards a generic method for plant viral disease diagnosis. *J virol Methods 2003*.
6. Song Y, Dai E, Wang J, Liu H, Zhai J,Chen C, Du Z, Guo Z, Yang R: Genotyping of hepatitis B virus HBV by oligonucleotides microarray. *Mol Cell Probes 2006*.
7. Perrin A, Duracher D, Perret M, Cleuziat P, Mandrand B:A combined oligonucleotide and protein microarray for the codetection of nucleic acids and antibodies associated with human immunodeficiency virus, hepatitis B virus and hepatitis C virus infections. *Anol Biochem 2003*.
8. Mezzasoma L, Crisanti A, Rossi R: Antigen Microarrays for Serodiagnosis of Infectious Diseases. *Clinical Chemistry 2002*.
9. Mohammed Z, Souna E, Anthony T: Principles of Bacterial Detection: Biosensors, Recognition Receptors and Microsystems.*Springer-Science and Business Media 2008*.
10. Wang D, Coscoy L, Zylberberg M, Avila PC, Boushey HA, Ganem D, DeRisi JL: Microarray-based detection and genotyping of viral pathogens. *Proc Natl Acad Sci USA 2002*.
11. Sergeev N, Distler M, Courtney S, Al-Khaldi SF, Volokhov D,Chizhikov V, Rasooly A: Multipathogen oligonucleotide microarray for environmental and biodefense applications. *Biosens Bioelectron 2005*.
12. Lemarchand K, Masson L, Brousseau R: Molecular biology and DNA microarray technology for microbial quality monitoring of water. *Crit Rev Microbiol 2004*.



13. Urism A, Fischer KF, Chiu CY, Kistler AL, Beck S, Wang D, DeRisi JL: E-Predict: a computational strategy for species identification based on observed DNA microarray hybridization patterns. *Genome Biol* 2005.
14. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: Basic Local Alignment Search tool. *J Mol Biol* 1990.
15. The R Project for Statistical Computing [<http://www.R-project.org>]
16. Michael Watson, Juliet Dukes, Abu-Bakr Abu-Median, Donald P King, Paul Britton: DetectiV: visualization, normalization and significance testing for pathogen-detection microarray data. *Genome Biol* 2007.
17. Smyth GK: Limma: linear models for microarray data. In *Bioinformatics and Computational Biology Solutions using R and Bioconductor* Edited by: Gentleman R, Carey V, Dudoit S, Irizarry R, Huber W. New York: Springer 2005
18. Gautier L, Cope L, Bolstad BM, Irizarry RA: affy - analysis of Affymetrix GeneChip data at the probe level. *Bioinformatics* 2004.
19. Tanya Barrett, Ron Edgar: Gene Expression Omnibus (GEO): Microarray data storage, submission, retrieval, and analysis. *Methods Enzymol* 2006.

Abstract

Microarray technology is the combination of different Sciences utilities such as Molecular Biology, Microelectronics, Microfluidics and Bioinformatics. By using this technology we can investigate, simultaneously, thousands of genetic or protein targets in a small system.

DNA microarrays consist of DNA microscopic points that are attached to a solid surface such as glass, plastic or silicon chip and formed as an array. The fixed pieces of DNA are considered as searchers. In an experiment, we can use thousands of searchers. Therefore, any microarray consists of the same number of genetic tests as the experiment performed on all of them in parallel. With this ability, arrays have speeded up the biological investigations. Microarray technology can be seen as a continued development of southern blotting.

However, the most important stage in this technology, analysis of data, requires reliable bioinformatics tools achieving high reliabilities.

Infectious diseases, from the beginning of human life, always have been with human beings and have caused major difficulties for them. One of the most important usages of microarray technology is the possibility of testing for the presence of thousands of micro-organism in environmental and clinical samples only in a single experiment, resulting in recognition of pathogens. Thereby, we can take an important step in curing the disease. As noted above, the most important stage in microarray technology is the data analysis. We present E-Predict algorithm and DetectiV package, for microarray based species identifications. We demonstrate the application of E-Predict and DetectiV to viral detection in a large, publicly available dataset and show that DetectiV performs better than E-Predict.

Keyword: 1-Microarray 2- Microarray Data 3- Microarray Data Analysis4- Infectious Diseases